

Simulation Studies on Some Nearest Neighbor Rules  
for Statistical Classification. <sup>(1)</sup>

By

David Aarons and Somesh Das Gupta

University of Minnesota  
Technical Report No. 303

November 1977

Simulation Studies on Some Nearest Neighbor Rules  
for Statistical Classification.<sup>(1)</sup>

by

David Aarons and Somesh DasGupta

University of Minnesota

<sup>(1)</sup> This research was supported by a grant from the Mathematics Division,  
U.S. Army, Research Triangle Park, N.C.; Grant DAAG 29-76-G-0038.

1. Introduction. The two-population classification problem is to identify a population  $\pi_0$  with one of two given populations  $\pi_1$  and  $\pi_2$  based on observations from these populations on a random vector  $X$ . We shall consider here  $X$  to be univariate. Let  $F_i$  be the c.d.f. of  $X$  in  $\pi_i$  ( $i = 0, 1, 2$ ). Thus our problem is to test  $H_1: F_0 = F_1$  vs.  $H_2: F_0 = F_2$ . In this paper we have considered some rules which are suggested in the literature when  $F_1, F_2$  are not known except that they are continuous. We have studied the performances of the following three rules by simulation.

Let  $X_0, X_{1i}$  ( $i = 1, \dots, n_1$ ),  $X_{2i}$  ( $i = 1, \dots, n_2$ ) be random observations on  $X$  from the populations  $\pi_0, \pi_1, \pi_2$ , respectively.

Rule I. 1-NN (nearest neighbor) Rule: Measure distances of  $X_0$  from  $X_{1i}$ 's and  $X_{2i}$ 's and based on these distances classify  $X_0$  into the population to which its nearest neighbor belongs.

Rule II. 1-RNN (rank nearest neighbor) Rule: Pool all the observations and order them.

(a) If  $X_0$  is the largest or the smallest observation classify  $X_0$  into the population of its nearest neighbor (based on ranks).

(b) If both the right-hand and the left-hand nearest neighbor of  $X_0$  (denoted by  $U_1$  and  $V_1$ ) belong to the same population, classify  $X_0$  into that population.

(c) If  $U_1$  and  $V_1$  belong to different populations classify  $X_0$  into  $\pi_1$  and  $\pi_2$  with probabilities  $1/2$  and  $1/2$ , respectively. (We call this case a "tie".)

Rule III. 2-RNN Rule: Apply the 1-RNN rule. If a tie occurs, delete the observations corresponding to  $U_1$  and  $V_1$  and apply the 1-RNN rule again on the remaining observations.

The first rule was suggested and studied by Fix and Hodges (1951, 1953). DasGupta and Lin (1977) proposed the RNN rules and obtained the asymptotic probabilities of misclassification as  $n_1, n_2 \rightarrow \infty$ . For a given rule  $\delta$ , let its PMC under  $F_0 = F_1$  be given by

$$\alpha(\delta) = \Pr[\delta \text{ classifies } X_0 \text{ into } \pi_2 | F_0 = F_1] .$$

Let  $\alpha_1^*, \alpha_2^*, \alpha_3^*$  be the asymptotic values of  $\alpha$  corresponding to the above rules 1, 2 and 3. Let  $f_i$  be the p.d.f. of  $F_i$  with respect to Lebesgue measure ( $i = 1, 2$ ) and  $p_i = \lim n_i / (n_1 + n_2)$  ( $i = 1, 2$ ) as  $\min(n_1, n_2) \rightarrow \infty$ . It was shown by Fix and Hodges (1951) and DasGupta and Lin (1977) that

$$\alpha_1^* = \alpha_2^* = \int_{-\infty}^{\infty} p_2 f_1(x) f_2(x) dx / \{p_1 f_1(x) + p_2 f_2(x)\}$$

$$\alpha_3^* = \alpha_2^* + \int_{-\infty}^{\infty} \frac{p_1 p_2 f_1(x) f_2(x) \cdot \{p_2 f_2(x) - p_1 f_1(x)\}}{\{p_1 f_1(x) + p_2 f_2(x)\}^3} f_1(x) dx .$$

In this paper we have studied the finite-sample performances of these rules by estimating  $\alpha$  based on samples from sets of two given populations.

2. The Experiment. Different steps of our simulation study are given below.

(i) Two known but different univariate distributions  $F_1$  and  $F_2$  are chosen.

(ii) Random samples of sizes  $n_1$  and  $n_2$  from  $F_1$  and  $F_2$ , respectively, are obtained; these samples are called training samples.

(iii) A random sample of size  $n_0$  from  $F_0 = F_1$  is obtained. We call this a test sample.

(iv) For each observation in the test sample a given classification rule  $\delta$  (one of the above three rules) is applied and let  $n_{02}$  be the number of the observations in the test sample which are classified by  $\delta$

into  $F_2$ . Let  $\hat{\alpha}(\delta) = n_{02}/n_0$  be the proportion of test samples misclassified into  $F_2$ .

(v) Steps (ii)-(iv) are repeated  $r$  times for new training and test samples keeping  $n_1$ ,  $n_2$  and  $n_0$  fixed.

(vi) The mean and the standard error of the mean based on  $r$  values of  $\hat{\alpha}(\delta)$  thus obtained are recorded.

(vii) Steps (ii)-(vi) are repeated for different values of  $n_1$ ,  $n_2$  and  $r$ .

(viii)  $F_2$  is characterized by a parameter  $\theta$ . For different values of  $\theta$  steps (i)-(vii) are repeated.

Our choices are given in the following table.

$F_1$	$F_2$	Parameters	$n_1=n_2$	$n_0$	$r$
$N(0,1)$	$N(\theta,1)$	$\theta=0, \pm 1, \pm 2, 3$	25 100	100 400	20 4
$N(0,1)$	$N(0,\theta)$	$\theta=2, 3, 1/2, 1/3$	25 100	100 400	20 4
$e^{-x}$ (density)	$\theta e^{-\theta x}$	$\theta=1, 2, 3, 4,$ $1/2, 1/3, 1/4, 1/8$	100	100	20
Cauchy (0,1)	Cauchy ( $\theta,1$ )	$\theta=0, \pm 1, \pm 2, \pm 3$	25 100	100 400	20 4

Samples are generated by a library subroutine available on the CDC 6400 at the University of Minnesota.

Note 1. In the following tables "Half" refers to taking one-half the number of ties to count as misclassified and "R-half" refers to resolving the ties by the use of uniform random number generator.

Note 2. In some of the following tables EPMC denotes an estimate of the asymptotic PMC ( $\alpha_1^* = \alpha_2^*$ ) of the 1-NN and 1-RNN rules. These are derived by the method of runs as suggested in Das Gupta and Lin (1977).

### 3. Tables

Table 3.1

Proportion of test sample misclassified into  $\pi_2$ .

$F_1 = N(0,1)$ ,  $F_2 = N(\theta,1)$ ;  $n_1 = n_2 = 25$ ,  $n_0 = 100$ ,  $r = 20$ .

Optimal (assuming  $\theta$  is known and for minimax rule) PMC is

$\Phi(-|\theta|/2)$ .

$\theta$ \ Rule	1NN		RNN		2-RNN		Opt. Exp't.
	MEAN	s.e.	MEAN	s.e.	MEAN	s.e.	
$\theta = 0$	.479	.017	Half .479 Rhalf .485	.013 .016	Half .479 Rhalf .484	.014 .015	.500
$\theta = 1$	.374	.018	Half .381 Rhalf .374	.014 .016	Half .343 Rhalf .340	.021 .021	.308
$\theta = -1$	.426	.020	Half .426 Rhalf .432	.014 .017	Half .421 Rhalf .425	.025 .024	.308
$\theta = 2$	.195	.018	Half .194 Rhalf .196	.018 .018	Half .165 Rhalf .164	.017 .018	.159
$\theta = -2$	.245	.020	Half .254 Rhalf .251	.018 .018	Half .258 Rhalf .255	.019 .018	.159
$\theta = 3$	.086	.012	Half .089 Rhalf .084	.012 .011	Half .062 Rhalf .061	.010 .009	.067
$\theta = -3$	.105	.013	Half .114 Rhalf .113	.012 .011	Half .119 Rhalf .118	.015 .015	.067

Table 3.2

Proportion of test sample misclassified into  $\pi_2$ .

$F_1 = N(0,1)$ ,  $F_2 = N(\theta,1)$ ;  $n_1 = n_2 = 100$ ,  $n_0 = 400$ ,  $r = 4$ .

$\theta$ \ Rule	1NN		RNN		EPMC	2-RNN		Opt. Exp't.
	MEAN	s.d.	MEAN	s.d.		MEAN	s.d.	PMC
$\theta = 0$	.490	.018	Half .482 .008 Rhalf .475 .006		.48	Half .509 .014 Rhalf .501 .016		.500
$\theta = 1$	.415	.010	Half .398 .014 Rhalf .404 .024		.36	Half .351 .009 Rhalf .358 .024		.308
$\theta = -1$	.402	.010	Half .394 .007 Rhalf .397 .007		.38	Half .347 .025 Rhalf .344 .024		.308
$\theta = 2$	.208	.010	Half .210 .010 Rhalf .208 .009		.22	Half .200 .011 Rhalf .199 .012		.159
$\theta = -2$	.209	.012	Half .213 .008 Rhalf .215 .009		.22	Half .197 .013 Rhalf .200 .014		.159
$\theta = 3$	.088	.011	Half .083 .009 Rhalf .082 .007		.10	Half .065 .005 Rhalf .066 .006		.007
$\theta = -3$	.104	.012	Half .101 .008 Rhalf .107 .013		.09	Half .088 .012 Rhalf .094 .014		.007

Table 3.3.

Proportion of test sample misclassified into  $\pi_2$ .

$F_1 = N(0,1)$ ,  $F_2 = N(0,\theta)$ ;  $n_1 = n_2 = 25$ ,  $n_0 = 100$ .  $r = 20$ .

$\theta$ \ Rule	1NN		RNN		2-RNN	
	MEAN	s.e.	MEAN	s.e.	MEAN	s.e.
$\theta = 2.0$	.375	.009	Half .394 Rhalf .393	.008 .010	Half .353 Rhalf .355	.014 .015
$\theta = 3.0$	.399	.014	Half .346 Rhalf .337	.013 .013	Half .293 Rhalf .295	.019 .018
$\theta = .5$	.417	.017	Half .438 Rhalf .337	.015 .018	Half .461 Rhalf .460	.020 .021
$\theta = 1/3$	.359	.022	Half .376 Rhalf .380	.018 .019	Half .393 Rhalf .391	.019 .019

Table 3.4

Proportion of test sample misclassified into  $\pi_2$ .

$F_1 = N(0,1)$ ,  $F_2 = N(0,\theta)$ ;  $n_1 = n_2 = 100$ ,  $n_0 = 400$ ,  $r = 4$ .

$\theta$ \ Rule	1NN		RNN		EPMC	2-RNN	
	MEAN	s.e.	MEAN	s.e.		MEAN	s.e.
$\theta = 2.0$	.435	.022	Half .424 Rhalf .426	.022 .025	.36	Half .395 Rhalf .396	.027 .028
$\theta = 3.0$	.333	.012	Half .338 Rhalf .336	.010 .011	.32	Half .295 Rhalf .296	.012 .011
$\theta = .5$	.397	.062	Half .405 Rhalf .407	.011 .013	.38	Half .409 Rhalf .408	.006 .005
$\theta = 1/3$	.339	.021	Half .352 Rhalf .354	.020 .021	.35	Half .360 Rhalf .361	.029 .030



Table 3.5

Proportion of test sample misclassified into  $\pi_2$ .

$$f_1(x) = e^{-x}, f_2(x) = \theta e^{-\theta x}; n_1 = n_2 = n_0 = 100, r = 4.$$

$\theta$ \ Rule	1NN		RNN		EPMC	2-RNN	
	MEAN	s.e.	MEAN	s.e.		MEAN	s.e.
$\theta = 1$	.508	.016	Half .509 Rhalf .523	.013 .016	.47	Half .503 Rhalf .517	.013 .015
$\theta = 2$	.442	.015	Half .434 Rhalf .438	.014 .017	.38	Half .442 Rhalf .444	.016 .016
$\theta = 3$	.402	.014	Half .388 Rhalf .387	.011 .011	.36	Half .394 Rhalf .387	.013 .014
$\theta = 4$	.335	.009	Half .330 Rhalf .336	.007 .008	.32	Half .327 Rhalf .330	.009 .009
$\theta = .5$	.453	.010	Half .453 Rhalf .458	.009 .013	.38	Half .430 Rhalf .430	.010 .014
$\theta = 1/3$	.410	.011	Half .395 Rhalf .386	.008 .009	.36	Half .346 Rhalf .335	.010 .010
$\theta = 1/4$	.354	.015	Half .364 Rhalf .372	.012 .013	.32	Half .290 Rhalf .292	.013 .013
$\theta = 1/8$	.247	.014	Half .248 Rhalf .259	.012 .014	.22	Half .181 Rhalf .185	.011 .010

Table 3.6

Proportion of test sample misclassified into  $\pi_2$ .

$F_1 = \text{Cauchy}(0,1)$ ,  $F_2 = \text{Cauchy}(\theta,1)$ ;  $n_1 = n_2 = 25$ ,  $n_0 = 100$ ,  $r = 20$ .

$\theta$ \ Rule	1NN		RNN		2-RNN	
	MEAN	s.e.	MEAN	s.e.	MEAN	s.e.
$\theta = 0$	.473	.018	Half .430 Rhalf .493	.015 .018	Half .488 Rhalf .505	.027 .029
$\theta = 1$	.406	.022	Half .418 Rhalf .408	.022 .025	Half .397 Rhalf .395	.031 .033
$\theta = -1$	.398	.016	Half .410 Rhalf .410	.012 .013	Half .389 Rhalf .385	.021 .022
$\theta = 2$	.288	.021	Half .297 Rhalf .288	.021 .021	Half .248 Rhalf .238	.027 .028
$\theta = -2$	.247	.012	Half .264 Rhalf .276	.012 .015	Half .248 Rhalf .252	.017 .019
$\theta = 3$	.161	.020	Half .168 Rhalf .161	.017 .018	Half .103 Rhalf .099	.017 .017
$\theta = -3$	.153	.015	Half .156 Rhalf .154	.013 .013	Half .130 Rhalf .125	.014 .014

Table 3.7

Proportion of test sample misclassified into  $\pi_2$ .

$F_1 = \text{Cauchy}(0,1)$ ,  $F_2 = \text{Cauchy}(\theta,1)$ ;  $n_1 = n_2 = 100$ ,  $n_0 = 400$ ,  $r = 4$ .

$\theta$ \ Rule	1NN		RNN			2-RNN		
	MEAN	s.e.		MEAN	s.e.		MEAN	s.e.
$\theta = 0$	.494	.015	Half	.514	.013	Half	.506	.017
			Rhalf	.529	.014	Rhalf	.512	.021
$\theta = 1$	.411	.010	Half	.426	.009	Half	.381	.018
			Rhalf	.446	.018	Rhalf	.390	.017
$\theta = -1$	.457	.029	Half	.446	.033	Half	.394	.028
			Rhalf	.454	.025	Rhalf	.393	.025
$\theta = 2$	.284	.007	Half	.278	.008	Half	.217	.033
			Rhalf	.283	.009	Rhalf	.219	.024
$\theta = -2$	.152	.016	Half	.318	.022	Half	.254	.014
			Rhalf	.321	.015	Rhalf	.257	.010
$\theta = 3$	.152	.016	Half	.154	.015	Half	.088	.018
			Rhalf	.417	.012	Rhalf	.087	.014
$\theta = -3$	.204	.034	Half	.199	.032	Half	.105	.011
			Rhalf	.198	.034	Rhalf	.103	.012

4. Concluding Remarks. For all the three rules considered, it seems that  $\hat{\alpha}_1$  has a definite tendency to decrease as  $\theta$  moves away (in either direction) from its value under  $F_1$ .

For small  $n_1 = n_2$  there is not any marked difference in performances of these three rules although the 2-RNN rule may be a bit better. However, for large  $n_1 = n_2$  the 2-RNN rule seem to have markedly better performance except for the cases  $N(0,1)$  vs.  $N(0,\theta)$ ,  $\theta < 1$ . This report is the first empirical study on the performances of 1NN and RNN rules, although a more detailed study especially on multi-stage RNN rules is called for.

References

1. Das Gupta, S. and Lin, H. E. (1977). Nearest neighbor rules for statistical classification based on ranks. Tech. Rep. 285, School of Statistics, University of Minnesota, Minneapolis, Minnesota.
2. Fix, E. and Hodges, J. L. (1951). Nonparametric discrimination: Consistency properties. U.S. Air Force School of Aviation Medicine. Report No. 4. Randolph Field, Texas.
3. Fix, E. and Hodges, J. L. (1953). Nonparametric discrimination. Small sample properties, Ibid., Report No. 11.